

DECENTRALIZED CONSENSUS ALGORITHM WITH DELAYED AND STOCHASTIC GRADIENTS

BENJAMIN SIRB AND XIAOJING YE

ABSTRACT. We analyze the convergence of decentralized consensus algorithm with delayed gradient information across the network. The nodes in the network privately hold parts of the objective function and collaboratively solve for the consensus optimal solution of the total objective while they can only communicate with their immediate neighbors. In real-world networks, it is often difficult and sometimes impossible to synchronize the nodes, and therefore they have to use stale gradient information during computations. We show that, as long as the random delays are bounded in expectation and a proper diminishing step size policy is employed, the iterates generated by decentralized gradient descent method converge to a consensual optimal solution. Convergence rates of both objective and consensus are derived. Numerical results on a number of synthetic problems and real-world seismic tomography datasets in decentralized sensor networks are presented to show the performance of the method.

1. INTRODUCTION

In this paper, we consider a decentralized consensus optimization problem arising from emerging technologies such as distributed machine learning [3, 10, 16, 19], sensor network [13, 29, 35], and smart grid [11, 21]. Given a network $G(V, E)$, $V = \{1, 2, \dots, m\}$ is the node (also called agent, processor, or sensor) set and $E \subset V \times V$ is the edge set. Two nodes i and j are called neighbors if $(i, j) \in E$. The communications between neighbor nodes are bidirectional, meaning that i and j can communicate with each other as long as $(i, j) \in E$.

In a decentralized sensor network G , individual nodes can acquire, store, and process data about large-sized objects. Each node i collects data and holds objective function $F_i(x; \xi_i)$ privately where $\xi_i \in \Theta$ is random with fixed but unknown probability distribution in domain Θ to model environmental fluctuations such as noise in data acquisition and/or inaccurate estimation of objective function or its gradient. Here $x \in X$ is the unknown (e.g., the seismic image) to be solved, where the domain $X \subset \mathbb{R}^n$ is compact and convex. Furthermore, we assume that $F_i(\cdot; \xi_i)$ is convex for all $\xi_i \in \Theta$ and $i \in V$, and we define $f_i(x) = \mathbb{E}_{\xi_i}[F_i(x; \xi_i)]$ which is convex with respect to $x \in X$. The goal of decentralized consensus optimization is to solve the minimization problem

$$(1) \quad \underset{x \in X}{\text{minimize}} f(x), \quad \text{where } f(x) := \sum_{i=1}^m f_i(x)$$

with the restrictions that $F_i(x; \xi_i)$, and hence $f_i(x)$, are accessible by node i only, and that nodes i and j can communicate only if $(i, j) \in E$ during the entire computation.

2000 *Mathematics Subject Classification.* 65K05, 90C25, 65Y05.

Key words and phrases. Decentralized consensus, delayed gradient, stochastic gradient, decentralized networks.

B. Sirb and X. Ye are with the Department of Mathematics and Statistics, Georgia State University, Atlanta, Georgia 30303, USA. E-mail: bsirb1@student.gsu.edu, xye@gsu.edu.

There are a number of practical issues that need to be taken into consideration in solving the real-world decentralized consensus optimization problem (1):

- The partial objective F_i (and f_i) is held privately by node i , and transferring F_i to a data fusion center is either infeasible or cost-ineffective due to data privacy, the large size of F_i , and/or limited bandwidth and communication power overhead of sensors. Therefore, the nodes can only communicate their own estimates of $x \in \mathbb{R}^n$ with their neighbors in each iteration of a decentralized consensus algorithm.
- Since it is often difficult and sometimes impossible for the nodes to be fully synchronized, they may not have access to the most up-to-date (stochastic) gradient information during computations. In this case, the node i has to use out-of-date (stochastic) gradient $\nabla F_i(x_i(t - \tau_i(t)); \xi_i(t - \tau_i(t)))$ where $x_i(t)$ is the estimate of x obtained by node i at iteration t , and $\tau_i(t)$ is the level of (possibly random) delay of the gradient information at t .
- The estimates $\{x_i(t)\}$ by the nodes should tend to be consensual as t increases, and the consensual value is a solution of problem (1). In this case, there is a guarantee of retrieving a good estimate of x from any surviving node in the network even if some nodes are sabotaged, lost, or run out of power during the computation process.

In this paper, we analyze a decentralized consensus algorithm which takes all the factors above into consideration in solving (1). We provide comprehensive convergence analysis of the algorithm, including the decay rates of objective function and disagreements between nodes, in terms of iteration number, level of delays, and network structure etc.

1.1. Related work. Distributed computing on networks is an emerging technology with extensive applications in modern machine learning [10, 16, 19], sensor networks [13, 29, 45, 46], and big data analysis [4, 30]. There are two types of scenarios in distributed computing: centralized and decentralized. In the centralized scenario, computations are carried out locally by worker (slave) nodes while computations of certain global variables must eventually be processed by designated master node or at a center of shared memory during each (outer) iteration. A major effort in this scenario has been devoted to update the global variable more effectively using an asynchronous setting in, for example, distributed centralized alternating direction method of multipliers (ADMM) [5, 7, 20, 40, 43]. In the decentralized scenario considered in this paper, the nodes privately hold parts of objective functions and can only communicate with neighbor nodes during computations. In many real-world applications, decentralized computing is particularly useful when a master-worker network setting is either infeasible or not economical, or the data acquisition and computation have to be carried out by individual nodes which then need to collaboratively solve the optimization problem. Decentralized networks are also more robust to node failure and can better address privacy concerns. For more discussions about motivations and advantages of decentralized computing, see, e.g., [15, 26, 28, 33, 37, 38] and references therein.

Decentralized consensus algorithms take the data distribution and communication restriction into consideration, so that they can be implemented at individual nodes in the network. In the *ideal synchronous case* of decentralized consensus where all the nodes are coordinated to finish computation and then start to exchange information with neighbors in each iteration, a number of developments have been made. A class of methods is to rewrite the consensus constraints for minimization problem (1) by introducing auxiliary variables between neighbor nodes (i.e., edges), and apply ADMM (possibly with linearization or preconditioning techniques) to derive an implementable decentralized consensus algorithm [6, 12, 14, 22, 34]. Most of these methods require each

node to solve a local optimization problem every iteration before communication, and reach a convergence rate of $O(1/T)$ in terms of outer iteration (communication) number T for general convex objective functions $\{f_i\}$. First-order methods based on decentralized gradient descent require less computational cost at individual nodes such that between two communications they only perform one step of a gradient descent-type update at the weighted average of previous iterates obtained from neighbors. In particular, Nesterov's optimal gradient scheme is employed in decentralized gradient descent with diminishing step sizes to achieve rate of $O(1/T)$ in [15], where an alternative gradient method that requires excessive communications in each inner iteration is also developed and can reach a theoretical convergence rate of $O(\log T/T^2)$, despite that it seems to work less efficiently in terms of communications than the former in practice. A correction technique is developed for decentralized gradient descent with convergence rate as $O(1/T)$ with constant step size in [33], which results in a saddle-point algorithm as pointed out in [23]. In [46], the authors combine Nesterov's gradient scheme and a multiplier-type auxiliary variable to obtain a fast optimality convergence rate of $O(1/T^2)$. Other first-order decentralized methods have also been developed recently, such dual averaging [8]. Additional constraints for primal variables in decentralized consensus optimization (1) are considered in [42].

In real-world decentralized computing, it is often difficult and sometimes impossible to coordinate all the nodes in the network such that their computation and communication are perfectly synchronized. One practical approach for such *asynchronous consensus* is using a broadcast scenario where in each (outer) iteration, one node in the network is assumed to wake up at random and broadcasts its value to neighbors (but does not hear them back). A number of algorithms for broadcast consensus are developed, for instance, in [2, 13, 24, 25]. Another important issue in the asynchronous setting is that nodes may have to use out-of-date (stale) gradient information during updates [27]. This delayed scenario in gradient descent is considered in a distributed but not decentralized setting in [1, 18, 36, 44]. In addition, analysis of stochastic gradient in distributed computing is also carried out in [1, 32]. In [9], linear convergence rate of optimality is derived for strongly convex objective functions with delays. Extending [1], a *fixed* delay at all nodes is considered in dual averaging [17] and gradient descent [39] in a decentralized setting, but they did not consider more practical and useful *random* delays, and there are no convergence rates on node consensus provided in these papers.

1.2. Contributions. The contribution of this paper is in three phases.

First, we consider a general decentralized consensus algorithm with randomly delayed and stochastic gradient (Section 2). In this case, the nodes do not need to be synchronized and they may only have access to stale gradient information. This renders stochastic gradients with random delays at different nodes in their gradient updates, which is suitable for many real-world decentralized computing applications.

Second, we provide a comprehensive convergence analysis of the proposed algorithm (Section 3). More precisely, we derive convergence rates for both the objective function (optimality) and disagreement (feasibility constraint of consensus), and show their dependency on the characteristics of the problem, such as Lipschitz constants of (stochastic) gradients and spectral gaps of the underlying network.

Third, we conduct a number of numerical experiments on synthetic and real datasets to validate the performance of the proposed algorithm (Section 4). In particular, we examine the convergence on synthetic decentralized least squares, robust least squares, and logistic regression problems. We also present the numerical results on the reconstruction of several seismic images in decentralized wireless sensor networks.

1.3. Notations and assumptions. In this paper, all vectors are column vectors unless otherwise noted. We denote by $x_i(t) \in \mathbb{R}^n$ the estimate of node i at iteration t , and $x(t) = (x_1(t), \dots, x_m(t))^T \in \mathbb{R}^{m \times n}$. We denote $\|x\| \equiv \|x\|_2$ if x is a vector and $\|x\| \equiv \|x\|_F$ if x is a matrix, which should be clear by the context. For any two vectors of same dimension, $\langle x, y \rangle$ denotes their inner product, and $\langle x, y \rangle_Q := \langle x, Qy \rangle$ for symmetric nonnegative definite matrix Q . For notation simplicity, we use $\langle x, y \rangle = \sum_{i=1}^m \langle x_i, y_i \rangle$ where x_i and y_i are the i -th row of the $m \times n$ matrices x and y respectively. Such matrix inner product is also generalized to $\langle x, y \rangle_Q$ for matrices x and y . In this paper, we set the domain $X := \{x \in \mathbb{R}^n : \|x\|_\infty \leq R\}$ for some $R > 0$, which can be thought of as the maximum pixel intensity in reconstructed images for instance. We further denote $\mathcal{X} := X^m \subset \mathbb{R}^{m \times n}$.

For each node i , we define $f_i(x) := \mathbb{E}_{\xi_i}[F_i(x; \xi_i)]$ as the expectation of objective function, and $g_i(t) := \nabla F_i(x(t); \xi_i(t))$ (here the gradient ∇ is taken with respect to x) is the stochastic gradient at $x_i(t)$ at node i . We let $\tau_i(t)$ be the delay of gradient at node i in iteration t , and $\tau(t) = (\tau_1(t), \dots, \tau_m(t))^T$. We write $f(x(t))$ in short for $\sum_{i=1}^m f_i(x_i(t)) \in \mathbb{R}$, $x(t - \tau(t))$ for $(x_1(t - \tau_1(t)), \dots, x_m(t - \tau_m(t)))^T \in \mathbb{R}^{m \times n}$, and $g(t - \tau(t))$ for $(g_1(t - \tau_1(t)), \dots, g_m(t - \tau_m(t)))^T \in \mathbb{R}^{m \times n}$. We assume f_i is continuously differentiable, ∇f_i has Lipschitz constant L_i , and denote $L := \max_{1 \leq i \leq m} L_i$. Let $x^* \in \mathbb{R}^n$ be a solution of (1). Since x^* is consensual, we denote $\mathbf{1}(x^*)^T$ simply by x^* in this paper which is clear by the context, for instance $f(x^*) = f(\mathbf{1}(x^*)^T) = \sum_{i=1}^m f_i(x^*)$. Furthermore, we let $y(T) := (1/T) \sum_{t=1}^T x(t+1)$ be the running average of $\{x(t+1) : 1 \leq t \leq T\}$, and $z(T) := (1/m) \sum_{i=1}^m y(T)$ be the consensus average of $y(T)$. We denote $J = (1/m) \mathbf{1} \mathbf{1}^T$, then $z(T) = Jy(T)$. Note that for all T , $z(T)$ is always consensual but $x(T), y(T)$ may not be.

Suppose $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable convex function, then for any $x, y \in \mathbb{R}^n$ we denote the Bregman distance (divergence) between x and y (order matters) by

$$(2) \quad D_h(x, y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle.$$

If in addition ∇h is L_h -Lipschitz continuous, then we can verify that, for any $x, y, z, w \in \mathbb{R}^n$, there is

$$(3) \quad \begin{aligned} \langle \nabla h(z) - \nabla h(w), x - y \rangle &= D_h(y, z) - D_h(x, z) - D_h(y, w) + D_h(x, w) \\ &\leq D_h(y, z) - D_h(x, z) + \frac{L_h}{2} \|x - w\|^2 \end{aligned}$$

where we used the facts that $D_h(y, w) \geq 0$ and $D_h(x, w) \leq \frac{L_h}{2} \|x - w\|^2$.

An important ingredient in decentralized gradient descent is the mixing matrix $W = [w_{ij}]$ in (4). For the algorithm to be implementable in practice, $w_{ij} > 0$ if and only if $(i, j) \in E$. In this paper, we assume that W is symmetric and $\sum_{j=1}^m w_{ij} = 1$ for all i , hence W is doubly stochastic, namely $W\mathbf{1} = \mathbf{1}$ and $\mathbf{1}^T W = \mathbf{1}^T$ where $\mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^m$. With the assumption that the network G is simple and connected, we know $\|W\|_2 = 1$ and eigenvalue 1 of W has multiplicity 1 by the Perron-Frobenius theorem. As a consequence, $Wx = x$ if and only if x is consensual, i.e., $x = c\mathbf{1}$ for some $c \in \mathbb{R}$. We further assume $W \succeq 0$ (otherwise use $\frac{1}{2}(I + W) \succeq 0$ since stochastic matrix W has spectral radius 1). Given a network G , there are different ways to design the mixing matrix W . For some optimal choices of W , see, e.g., [31, 41].

Now we make several mild assumptions that are necessary in our convergence analysis.

- (1) The network $G(V, E)$ is undirected, simple, and connected.
- (2) The stochastic gradient satisfies $\mathbb{E}_{\xi_i}[\nabla F_i(x; \xi_i)] = \nabla f_i(x)$ for all i and x . Moreover, for all i , ξ , and x , $\|\nabla f_i\|$ and $\mathbb{E}_{\xi_i}[\|\nabla F_i(x; \xi_i)\|^2] \leq G^2$ for some $G > 0$, and $\mathbb{E}_{\xi_i}[\|\nabla F_i(x; \xi_i) - \nabla f_i(x)\|^2] \leq \sigma^2$ for some $\sigma > 0$.

- (3) The delays $\tau_i(t)$ may follow different distributions at different nodes, but their second moments are assumed to be uniformly bounded, i.e., there exists $B > 0$ such that $\mathbb{E}[|\tau_i(t)|^2] \leq B^2$ for all $i = 1, \dots, m$ and iteration t . For each node i , we assume each update happens once, i.e., $t \mapsto t - \tau_i(t)$ is strictly increasing as t increases.

It is worth pointing out that these assumptions are rather standard and easy to satisfy in practice. For instance, the boundedness of ∇f_i is a consequence of the compactness of domain X and the Lipschitz continuity of ∇f_i . The assumption on random delays in a distributed system is also used in [1]. We further assume that the stochastic error ξ_i and the random delay τ_i are independent.

2. ALGORITHM

Taking the delayed stochastic gradient and the constraint that nodes can only communicate with immediate neighbors, we propose the following decentralized delayed stochastic gradient descent method for solving (1). Starting from an initial guess $\{x_i(0) : i = 1, \dots, m\}$, each node i performs the following updates iteratively:

$$(4) \quad x_i(t+1) = \Pi_X \left[\sum_{j=1}^m w_{ij} x_j(t) - \alpha(t) g_i(t - \tau_i(t)) \right].$$

Namely, in each iteration t , the nodes exchange their most recent $x_i(t)$ with their neighbors. Then each node takes weighted average of the received local copies using weights w_{ij} and performs a gradient descent type update using a stochastic gradient $g_i(t - \tau_i(t))$ with delay $\tau_i(t)$ and step size $\alpha(t)$, and projects the result onto X .

Following the matrix notation in Section 1.3, the iteration (4) can be written as

$$(5) \quad x(t+1) = \Pi_X [Wx(t) - \alpha(t)g(t - \tau(t))].$$

Here the projection Π_X is accomplished by each node projecting to X due to the definition of X in Section 1.3, which does not require any coordination between nodes. Note that the update (5) is also equivalent to

$$(6) \quad x(t+1) = \operatorname{argmin}_{x \in X} \left\{ \langle g(t - \tau(t)), x \rangle + \frac{1}{2\alpha(t)} \|x - Wx(t)\|^2 \right\}.$$

In this paper, we may refer to the proposed decentralized delayed stochastic gradient descent algorithm by any of (4), (5), and (6) since they are equivalent.

3. CONVERGENCE ANALYSIS

In this section, we provide a comprehensive convergence analysis of the proposed algorithm (6) by employing a proper step size policy. In particular, we derive convergence rates for the objective function and disagreement in that order.

Lemma 1. *Let $\{x(t)\}$ be the iterates generated by Algorithm (5), then the following inequality holds for all $T \geq 1$:*

$$(7) \quad \sum_{t=1}^T \langle \nabla f(x(t)) - \nabla f(x(t - \tau(t))), x(t+1) - x^* \rangle \\ \leq 2mnLR^2(1 + 2B^2) + \frac{L}{2}(B+1)^2 \sum_{t=1}^T \|x(t+1) - x(t)\|^2.$$

Proof. We first observe that

$$\begin{aligned}
 & \langle \nabla f(x(t)) - \nabla f(x(t - \tau(t))), x(t+1) - x^* \rangle \\
 (8) \quad &= \sum_{i=1}^m \langle \nabla f_i(x_i(t)) - \nabla f_i(x_i(t - \tau_i(t))), x_i(t+1) - x^* \rangle \\
 &\leq \sum_{i=1}^m \left[D_{f_i}(x^*, x_i(t)) - D_{f_i}(x^*, x_i(t - \tau_i(t))) + \frac{L}{2} \|x_i(t+1) - x_i(t - \tau_i(t))\|^2 \right]
 \end{aligned}$$

where we applied (3) to get the inequality. We further note that the convexity of $\|\cdot\|^2$ implies

$$(9) \quad \|x_i(t+1) - x_i(t - \tau_i(t))\|^2 \leq (\tau_i(t) + 1) \sum_{s=0}^{\tau_i(t)} \|x_i(t-s+1) - x_i(t-s)\|^2.$$

Combining (8) and (9), and taking the sum of t from 1 to T , we obtain

$$\begin{aligned}
 & \sum_{t=1}^T \langle \nabla f(x(t)) - \nabla f(x(t - \tau(t))), x(t+1) - x^* \rangle \\
 (10) \quad &\leq \sum_{i=1}^m \left[\sum_{t=1}^T (D_{f_i}(x^*, x_i(t)) - D_{f_i}(x^*, x_i(t - \tau_i(t)))) \right. \\
 &\quad \left. + \frac{L}{2} \sum_{t=1}^T (\tau_i(t) + 1) \sum_{s=0}^{\tau_i(t)} \|x_i(t-s+1) - x_i(t-s)\|^2 \right]
 \end{aligned}$$

For each i , the sum of D_{f_i} terms for t from 1 to T above leaves only those not received by the gradient procedure within T iterations, namely

$$(11) \quad \sum_{t=1}^T (D_{f_i}(x^*, x_i(t)) - D_{f_i}(x^*, x_i(t - \tau_i(t)))) = \sum_{t \in \mathcal{S}_i(T)} D_{f_i}(x^*, x_i(t))$$

where $\mathcal{S}_i(T) := \{1 \leq t \leq T : t > T - \tau_i(T)\}$. Then by Chebyshev's inequality, we can bound the expected cardinality of $\mathcal{S}_i(T)$ by

$$(12) \quad \mathbb{E}[|\mathcal{S}_i(T)|] = \sum_{t=1}^T \mathbb{P}(\tau_i(T) > T - t) \leq 1 + \sum_{t=1}^{T-1} \frac{B^2}{(T-t)^2} \leq 1 + 2B^2.$$

where we used the fact that $\sum_{t=1}^{T-1} \frac{1}{(T-t)^2} = \sum_{t=1}^{T-1} \frac{1}{t^2} \leq 2 - \frac{1}{T-1} \leq 2$. Combining (11) and (12), and using the fact that $D_{f_i}(x^*, x_i(t)) \leq 2nLR^2$ for all t and i , we obtain,

$$(13) \quad \sum_{i=1}^m \sum_{t=1}^T (D_{f_i}(x^*, x_i(t)) - D_{f_i}(x^*, x_i(t - \tau_i(t)))) \leq 2mnLR^2(1 + 2B^2).$$

For each i , the second sum for t from 1 to T on the right side of (10) yields

$$(14) \quad \sum_{t=1}^T (\tau_i(t) + 1) \sum_{s=0}^{\tau_i(t)} \|x_i(t-s+1) - x_i(t-s)\|^2 \leq \sum_{t=1}^T N_i(t, T) \|x_i(t+1) - x_i(t)\|^2$$

where the coefficient $N_i(t, T)$ is defined by

$$(15) \quad N_i(t, T) := \sum_{\{t \leq s \leq T: 0 \leq s - \tau_i(s) \leq t\}} (\tau_i(s) + 1)$$

Therefore, we have for each i that

$$(16) \quad \begin{aligned} \mathbb{E}[N_i(t, T)] &= \mathbb{E} \left[\sum_{\{t \leq s \leq T: 0 \leq s - \tau_i(s) \leq t\}} (\tau_i(s) + 1) \right] = \sum_{s=t}^T \sum_{k=s-t}^s (k+1) \mathbb{P}(\tau_i(s) = k) \\ &\leq \sum_{k=0}^T (k+1)^2 \mathbb{P}(\tau_i(s) = k) \leq \mathbb{E}[|\tau_i(s) + 1|^2] \leq (B+1)^2 \end{aligned}$$

where the first inequality is obtained by listing each possible value of k in the double sum, and upper bounding its occurrence by $(k+1)$, and the last inequality is due to $\mathbb{E}[|\tau_i(s)|] \leq \sqrt{\mathbb{E}[|\tau_i(s)|^2]} = B$. Therefore, (14) can be bounded by

$$(17) \quad \sum_{i=1}^m \sum_{t=1}^T (\tau_i(t) + 1) \sum_{s=0}^{\tau_i(t)} \|x_i(t-s+1) - x_i(t-s)\|^2 \leq (B+1)^2 \sum_{t=1}^T \|x(t+1) - x(t)\|^2$$

Applying (13) and (17) to (10) completes the proof. \square

Theorem 2. *Let $\{x(t)\}$ be the iterates generated by Algorithm (5) with $\alpha(t) = [2(L + \eta(t))]^{-1}$ where $\eta(t)$ is a nondecreasing function of t , then*

$$(18) \quad \begin{aligned} \mathbb{E}[f(y(T)) - f(x^*)] &\leq \frac{2mnR^2[4L + 2\eta(1) + 2\eta(T) + L(1 + 2B^2)]}{T} + \frac{2m\sigma^2}{T} \sum_{t=1}^T \frac{1}{\eta(t)} \\ &\quad + \frac{L(B+1)^2}{2T} \sum_{t=1}^T \mathbb{E}[\|x(t+1) - x(t)\|^2] \end{aligned}$$

where $y(T) = (1/T) \sum_{t=1}^T x(t+1)$ is the running average of $\{x(t)\}$.

Proof. We first note that there is

$$\begin{aligned}
(19) \quad & f(x(t+1)) - f(x^*) = \sum_{i=1}^m (f_i(x_i(t+1)) - f_i(x^*)) \\
&= \sum_{i=1}^m [f_i(x_i(t+1)) - f_i(x_i(t)) + f_i(x_i(t)) - f_i(x^*)] \\
&\leq \sum_{i=1}^m \left[\langle \nabla f_i(x_i(t)), x_i(t+1) - x_i(t) \rangle + \frac{L_i}{2} \|x_i(t+1) - x_i(t)\|^2 \right. \\
&\quad \left. + \langle \nabla f_i(x_i(t)), x_i(t) - x^* \rangle \right] \\
&\leq \sum_{i=1}^m \left[\langle \nabla f_i(x_i(t)), x_i(t+1) - x^* \rangle + \frac{L_i}{2} \|x_i(t+1) - x_i(t)\|^2 \right] \\
&\leq \langle \nabla f(x(t)), x(t+1) - x^* \rangle + \frac{L}{2} \|x(t+1) - x(t)\|^2 \\
&\leq \langle g(t - \tau(t)), x(t+1) - x^* \rangle + \langle \nabla f(x(t)) - g(t - \tau(t)), x(t+1) - x^* \rangle \\
&\quad + \frac{L}{2} \|x(t+1) - x(t)\|^2
\end{aligned}$$

where we used the L_i -Lipschitz continuity of ∇f_i and convexity of f_i to obtain the first inequality. Note that $x(t+1)$ is obtained by (6) as

$$\begin{aligned}
(20) \quad & x(t+1) = \operatorname{argmin}_{x \in \mathcal{X}} \left\{ \langle g(t - \tau(t)), x \rangle + \frac{1}{2\alpha(t)} \|x - Wx(t)\|^2 \right\} \\
&= \operatorname{argmin}_{x \in \mathcal{X}} \left\{ \left\langle g(t - \tau(t)) + \frac{1}{\alpha(t)} (I - W)x(t), x \right\rangle + \frac{1}{2\alpha(t)} \|x - x(t)\|^2 \right\}
\end{aligned}$$

Therefore, the optimality of $x(t+1)$ in (6) implies that

$$\begin{aligned}
(21) \quad & \langle g(t - \tau(t)), x(t+1) - x^* \rangle \\
&\leq -\frac{1}{\alpha(t)} \langle (I - W)x(t), x(t+1) - x^* \rangle \\
&\quad + \frac{1}{2\alpha(t)} [\|x^* - x(t)\|^2 - \|x(t+1) - x(t)\|^2 - \|x^* - x(t+1)\|^2].
\end{aligned}$$

Furthermore, we note that $\mathbf{1} \in \operatorname{Null}(I - W)$ and x^* is consensual, hence we have

$$\begin{aligned}
(22) \quad & -\frac{1}{\alpha(t)} \langle (I - W)x(t), x(t+1) - x^* \rangle \\
&= -\frac{1}{\alpha(t)} \langle (I - W)(x(t) - x^*), x(t+1) - x^* \rangle \\
&= \frac{1}{2\alpha(t)} (\|x(t) - x(t+1)\|_{I-W}^2 - \|x(t) - x^*\|_{I-W}^2 - \|x(t+1) - x^*\|_{I-W}^2) \\
&\leq \frac{1}{4\alpha(t)} \|x(t) - x(t+1)\|_{I-W}^2
\end{aligned}$$

where we have used the fact that

$$\|x(t) - x(t+1)\|_{I-W}^2 \leq 2(\|x(t) - x^*\|_{I-W}^2 + \|x(t+1) - x^*\|_{I-W}^2)$$

to obtain the inequality above. We also have that

$$\|x(t) - x(t+1)\|_{I-W}^2 \leq \|x(t) - x(t+1)\|^2$$

with which we can further bound (22) as

$$-\frac{1}{\alpha(t)} \langle (I-W)x(t), x(t+1) - x^* \rangle \leq \frac{1}{4\alpha(t)} \|x(t) - x(t+1)\|^2.$$

Now applying the inequality above and (21) to (19), and taking sum of t from 1 to T , we get

$$\begin{aligned} & \sum_{t=1}^T f(x(t+1)) - Tf(x^*) \\ (23) \quad & \leq \sum_{t=1}^T \left[\frac{1}{2\alpha(t)} (\|x(t) - x^*\|^2 - \|x(t+1) - x^*\|^2) + \left(\frac{L}{2} - \frac{1}{4\alpha(t)} \right) \|x(t) - x(t+1)\|^2 \right] \\ & \quad + \sum_{t=1}^T \langle \nabla f(x(t)) - g(t - \tau(t)), x(t+1) - x^* \rangle. \end{aligned}$$

For the last term on the right hand side of (23), we have

$$\begin{aligned} & \sum_{t=1}^T \langle \nabla f(x(t)) - g(t - \tau(t)), x(t+1) - x^* \rangle \\ & = \sum_{t=1}^T \langle \nabla f(x(t)) - \nabla f(x(t - \tau(t))), x(t+1) - x^* \rangle \\ (24) \quad & \quad + \sum_{t=1}^T \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x^* \rangle \\ & \leq 2mnLR^2(1 + 2B^2) + \sum_{t=1}^T \frac{L}{2} (B+1)^2 \|x(t+1) - x(t)\|^2 \\ & \quad + \sum_{t=1}^T \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x^* \rangle \end{aligned}$$

where we applied the Lemma 1 to obtain the inequality.

Note that the running average $y(T) = (1/T) \sum_{t=1}^T x(t+1)$ satisfies $f(y(T)) \leq \sum_{t=1}^T f(x(t+1))$ due to the convexity of all f_i . Therefore, together with (23) and (24) and the definition of $\alpha(t)$, we have

$$\begin{aligned} & T(f(y(T)) - f(x^*)) \\ (25) \quad & \leq \sum_{t=1}^T \left[\frac{1}{2\alpha(t)} (\|x(t) - x^*\|^2 - \|x(t+1) - x^*\|^2) + \frac{L(B+1)^2 - \eta(t)}{2} \|x(t) - x(t+1)\|^2 \right] \\ & \quad + 2mnLR^2(1 + 2B^2) + \sum_{t=1}^T \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x^* \rangle. \end{aligned}$$

Now, by taking expectation on both sides of (25), we obtain

$$(26) \quad \begin{aligned} T \mathbb{E}[f(y(T)) - f(x^*)] &\leq \sum_{t=1}^T \left[\frac{1}{2\alpha(t)} (e(t) - e(t+1)) + \frac{L(B+1)^2 - \eta(t)}{2} \mathbb{E}[\|x(t) - x(t+1)\|^2] \right] \\ &\quad + 2mnLR^2(1 + 2B^2) + \sum_{t=1}^T \mathbb{E} \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x^* \rangle \end{aligned}$$

where we denoted $e(t) := \mathbb{E}[\|x(t) - x^*\|^2]$ for notation simplicity.

Now we work on the last sum of inner products on the right side of (26). First we observe that

$$(27) \quad \begin{aligned} &\mathbb{E} \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x^* \rangle \\ &= \mathbb{E} \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t) - x^* \rangle \\ &\quad + \mathbb{E} \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x(t) \rangle. \end{aligned}$$

Note that $g(t - \tau(t))$ is used to calculate $x(t+1)$, and hence its stochastic error $g(t - \tau(t)) - \nabla f(x(t - \tau(t)))$ is independent of $x(t)$. Therefore, we have

$$(28) \quad \mathbb{E} \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t) - x^* \rangle = 0.$$

Furthermore, by Young's inequality, we have

$$(29) \quad \begin{aligned} &\mathbb{E} \langle \nabla f(x(t - \tau(t))) - g(t - \tau(t)), x(t+1) - x(t) \rangle \\ &\leq \frac{2}{\eta(t)} \mathbb{E}[\|\nabla f(x(t - \tau(t))) - g(t - \tau(t))\|^2] + \frac{\eta(t)}{2} \mathbb{E}[\|x(t+1) - x(t)\|^2] \\ &\leq \frac{2m\sigma^2}{\eta(t)} + \frac{\eta(t)}{2} \mathbb{E}[\|x(t+1) - x(t)\|^2] \end{aligned}$$

where we used the fact that $\mathbb{E}[\|\nabla f(x(t - \tau(t))) - g(t - \tau(t))\|^2] \leq m\sigma^2$ for all t . Now applying (27), (28) and (29) in (26), we have

$$(30) \quad \begin{aligned} &T \mathbb{E}[f(y(T)) - f(x^*)] \\ &\leq \sum_{t=1}^T \frac{1}{2\alpha(t)} (e(t) - e(t+1)) + 2mnLR^2(1 + 2B^2) + \sum_{t=1}^T \frac{2m\sigma^2}{\eta(t)} \\ &\quad + \frac{L(B+1)^2}{2} \sum_{t=1}^T \mathbb{E}[\|x(t+1) - x(t)\|^2] \\ &\leq \frac{e(1)}{2\alpha(1)} + \sum_{t=2}^T \frac{e(t)}{2} \left(\frac{1}{\alpha(t)} - \frac{1}{\alpha(t-1)} \right) + 2mnLR^2(1 + 2B^2) + \sum_{t=1}^T \frac{2m\sigma^2}{\eta(t)} \\ &\quad + \frac{L(B+1)^2}{2} \sum_{t=1}^T \mathbb{E}[\|x(t+1) - x(t)\|^2]. \end{aligned}$$

Note that $\alpha(t)$ is nonincreasing, therefore $\frac{1}{\alpha(t)} - \frac{1}{\alpha(t-1)} \geq 0$ and hence

$$(31) \quad \sum_{t=2}^T \frac{e(t)}{2} \left(\frac{1}{\alpha(t)} - \frac{1}{\alpha(t-1)} \right) \leq 2mnR^2 \sum_{t=2}^T \left(\frac{1}{\alpha(t)} - \frac{1}{\alpha(t-1)} \right) \leq \frac{2mnR^2}{\alpha(T)}$$

where we used the fact that $e(t) = \mathbb{E}[\|x(t) - x^*\|^2] \leq 4mnR^2$ for all t . Applying (31) to (30) yields (18). \square

We have shown that the running average $y(T)$ makes the objective function decay as in (18). However, an important feature in decentralized computing is that $x_i(t)$ tend to be consensual. Now we prove that the consensus can be achieved by the proposed algorithm (5), and we derive the convergence rate for the employed step size policy.

Lemma 3. *For any $x \in \mathbb{R}^{m \times n}$, its projection onto \mathcal{X} yields nonincreasing disagreement. That is*

$$(32) \quad \|(I - J) \Pi_{\mathcal{X}}(x)\|^2 \leq \|(I - J)x\|^2$$

Proof. See Appendix A. \square

Lemma 4. *Let $c_1 \geq 0$ and $c_2 > 0$, and define $\alpha(t) = 1/(c_1 + c_2\sqrt{t})$. Then for any $\lambda \in (0, 1)$ there is*

$$(33) \quad \sum_{s=0}^{t-1} \alpha(s) \lambda^{t-s-1} \leq \frac{\sqrt{\pi} \lambda^{-2}}{c_2 \sqrt{t} \log(\lambda^{-1})} = O\left(\frac{1}{\sqrt{t}}\right)$$

for all $t = 1, 2, \dots$

Proof. See Appendix B. \square

Theorem 5. *Let $\{x(t)\}$ be the iterates generated by Algorithm (6) with $\alpha(t) = [2(L + \eta\sqrt{t})]^{-1}$ for $\eta > 0$, and $\lambda = \|W - J\|$. Then λ is the second largest eigenvalue of W and hence $\lambda \in (0, 1)$. Moreover, the disagreement of $x(t)$ is bounded above by*

$$(34) \quad \|(I - J)x(t)\| \leq \sqrt{m}G \sum_{s=0}^{t-1} \alpha(s) \lambda^{t-s-1} \leq \frac{\sqrt{\pi m}G \lambda^{-2}}{\eta \sqrt{t} \log(\lambda^{-1})}$$

and the disagreement of running average $y(T) = (1/m) \sum_{t=1}^T x(t+1)$ is bounded above by

$$(35) \quad \|(I - J)y(T)\| \leq \frac{2\sqrt{\pi m}G \lambda^{-2}}{\eta \sqrt{T} \log(\lambda^{-1})} = O\left(\frac{1}{\sqrt{T}}\right).$$

Proof. We prove this bound by induction. It is trivial to show the bound for $t = 1$. Assume (34) is true for t , then we have

$$(36) \quad \begin{aligned} \|(I - J)x(t+1)\| &= \|(I - J) \Pi_{\mathcal{X}}(Wx(t) - \alpha(t)g(t - \tau(t)))\| \\ &\leq \|(I - J)(Wx(t) - \alpha(t)g(t - \tau(t)))\| \\ &\leq \|(I - J)Wx(t)\| + \alpha(t)\|(I - J)g(t - \tau(t))\| \\ &\leq \|(I - J)Wx(t)\| + \alpha(t)\sqrt{m}G \end{aligned}$$

where we used Lemma 3 in the first inequality, and $\|I - J\| \leq 1$ and $\|g_i(t - \tau_i(t))\| \leq G$ in the last inequality. Noting that $J^2 = J$ and $JW = WJ = J$, we have

$$(W - J)(I - J) = (I - J)W.$$

Therefore, we obtain

$$\begin{aligned}
(37) \quad \|(I - J)x(t + 1)\| &\leq \|(I - J)Wx(t)\| + \alpha(t)\sqrt{m}G \\
&= \|(W - J)(I - J)x(t)\| + \alpha(t)\sqrt{m}G \\
&\leq \|(W - J)\| \|(I - J)x(t)\| + \alpha(t)\sqrt{m}G \\
&\leq \lambda\sqrt{m}G \sum_{s=0}^{t-1} \alpha(s)\lambda^{t-s-1} + \alpha(t)\sqrt{m}G \\
&= \sqrt{m}G \sum_{s=0}^t \alpha(s)\lambda^{t-s}
\end{aligned}$$

where we used the induction assumption for t in the last inequality. Applying Lemma 4 yields the second inequality in (34). By convexity of $\|\cdot\|$ and definition of $y(T)$, we have

$$(38) \quad \|(I - J)y(T)\| \leq \frac{1}{T} \sum_{t=1}^T \|(I - J)x(t + 1)\| \leq \sum_{t=1}^T \frac{\sqrt{\pi m}G\lambda^{-2}}{\eta\sqrt{t}\log(\lambda^{-1})} \leq \frac{2\sqrt{\pi m}G\lambda^{-2}}{\eta\sqrt{T}\log(\lambda^{-1})}$$

by applying (34) and using $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$. \square

Corollary 6. *Let $\{x(t)\}$ be the iterates generated by Algorithm (6) with the settings of $\alpha(t)$, λ , and η same as in Theorem 5. Then there is*

$$(39) \quad \|x(t + 1) - x(t)\|^2 \leq \frac{2G^2}{\eta^2 t} \left[\frac{\pi m \lambda^{-4}}{\log^2(\lambda^{-1})} + \frac{1}{4} \right]$$

Proof. According to the update (6) or equivalently (5), we have

$$\begin{aligned}
(40) \quad \|x(t + 1) - x(t)\|^2 &= \|\Pi_{\mathcal{X}}(Wx(t) - \alpha(t)g(t - \tau(t))) - x(t)\|^2 \\
&\leq \|(I - W)x(t) + \alpha(t)g(t - \tau(t))\|^2 \\
&\leq 2(\|(I - W)x(t)\|^2 + \|\alpha(t)g(t - \tau(t))\|^2)
\end{aligned}$$

where we used the facts that $x(t) \in \mathcal{X}$ and that projection $\Pi_{\mathcal{X}}$ is nonexpansive. Note that $WJ = J$ and hence $I - W = (I - W)(I - J)$, we have

$$(41) \quad \|(I - W)x(t)\|^2 = \|(I - W)(I - J)x(t)\|^2 \leq \|(I - J)x(t)\|^2 \leq \frac{\pi m G^2 \lambda^{-4}}{\eta^2 t \log^2(\lambda^{-1})}$$

where we used the fact that $\|I - W\| \leq 1$ in the first inequality and applied Theorem 5 to obtain the second inequality.

On the other hand, we have by the definition of $\alpha(t)$ that

$$(42) \quad \|\alpha(t)g(t - \tau(t))\|^2 \leq (\alpha(t))^2 G^2 = \frac{G^2}{4(L + \eta\sqrt{t})^2} \leq \frac{G^2}{4\eta^2 t}$$

Applying (41) and (42) to (40) yields (39). \square

Theorem 7. *Let $x(t)$ be generated by Algorithm (4) with $\alpha(t) = [2(L + \eta\sqrt{t})]^{-1}$ for some $\eta > 0$. Let $y(T) = (1/T) \sum_{t=1}^T x(t + 1)$ be the running average of $x(t)$ and $z(T) = Jy(T) = (1/m) \sum_{i=1}^m y_i(T)$*

be the consensus average of $y(T)$, then

$$\begin{aligned}
 & \mathbb{E}[f(z(T))] - f(x^*) \\
 (43) \quad & \leq \frac{2\sqrt{\pi}mG^2}{\eta\lambda^2\sqrt{T}\log(\lambda^{-1})} + \frac{2mnR^2(4L + 2\eta + L(1 + 2B^2))}{T} + \frac{2mnR^2\eta}{\sqrt{T}} + \frac{4m\sigma^2}{\eta\sqrt{T}} \\
 & \quad + \frac{2L(B+1)^2G^2(1 + \log T)}{\eta^2T} \left[\frac{\pi m\lambda^{-4}}{\log^2(\lambda^{-1})} + \frac{1}{4} \right].
 \end{aligned}$$

Proof. We first bound the difference between the function values at the running average $y(T)$ and the consensus average $z(T) = Jy(T)$:

$$\begin{aligned}
 & |f(y(T)) - f(z(T))| = \left| \sum_{i=1}^m (f_i(y_i(T)) - f_i(z(T))) \right| \leq \sum_{i=1}^m |\langle \nabla f_i(z(T)), y_i(T) - z(T) \rangle| \\
 (44) \quad & \leq G \sum_{i=1}^m \|y_i(T) - z(T)\| \leq \sqrt{m}G \|(I - J)y(T)\| \leq \frac{\sqrt{m}G}{T} \sum_{t=1}^T \|(I - J)x(t+1)\| \\
 & \leq \frac{\sqrt{\pi}mG^2}{\eta\lambda^2T\log(\lambda^{-1})} \sum_{t=1}^T \frac{1}{\sqrt{t}} \leq \frac{2\sqrt{\pi}mG^2}{\eta\lambda^2\sqrt{T}\log(\lambda^{-1})}
 \end{aligned}$$

where we used convexity of f_i in the first inequality, $\|\nabla f_i\| \leq G$ in the second inequality, the convexity of $\|\cdot\|$ in the fourth inequality, and Theorem 5 to get the fifth inequality. Note that Theorem 2 implies

$$\begin{aligned}
 (45) \quad \mathbb{E}[f(y_T) - f(x^*)] & \leq \frac{2mnR^2(2L + \eta + L(1 + 2B^2))}{T} + \frac{2mnR^2\eta}{\sqrt{T}} + \frac{4m\sigma^2}{\eta\sqrt{T}} \\
 & \quad + \frac{L(B+1)^2}{2T} \sum_{t=1}^T \mathbb{E}[\|x(t+1) - x(t)\|^2],
 \end{aligned}$$

and the last term on the right hand side can be bounded by using Corollary 6:

$$\begin{aligned}
 (46) \quad \frac{L(B+1)^2}{2T} \sum_{t=1}^T \mathbb{E}[\|x(t+1) - x(t)\|^2] & \leq \frac{2L(B+1)^2G^2}{\eta^2T} \left[\frac{\pi m\lambda^{-4}}{\log^2(\lambda^{-1})} + \frac{1}{4} \right] \sum_{t=1}^T \frac{1}{t} \\
 & \leq \frac{2L(B+1)^2G^2(1 + \log T)}{\eta^2T} \left[\frac{\pi m\lambda^{-4}}{\log^2(\lambda^{-1})} + \frac{1}{4} \right]
 \end{aligned}$$

where we used the fact that $\sum_{t=1}^T (1/t) \leq 1 + \log T$. Therefore, we obtain

$$\begin{aligned}
 & \mathbb{E}[f(z(T))] - f(x^*) \leq \mathbb{E}[|f(z(T)) - f(y(T))|] + \mathbb{E}[f(y_T) - f(x^*)] \\
 (47) \quad & \leq \frac{2\sqrt{\pi}mG^2}{\eta\lambda^2\sqrt{T}\log(\lambda^{-1})} + \frac{2mnR^2(2L + \eta + L(1 + 2B^2))}{T} + \frac{2mnR^2\eta}{\sqrt{T}} + \frac{4m\sigma^2}{\eta\sqrt{T}} \\
 & \quad + \frac{2L(B+1)^2G^2(1 + \log T)}{\eta^2T} \left[\frac{\pi m\lambda^{-4}}{\log^2(\lambda^{-1})} + \frac{1}{4} \right]
 \end{aligned}$$

On the other hand, $z(T)$ is consensus, so $f(z(T)) \geq f(x^*)$ since x^* is a solution of (1). This completes the proof. \square

4. NUMERICAL EXPERIMENTS

In this section, we test algorithm (4) on decentralized consensus optimization problem (1) with delayed stochastic gradients using a number of synthetic and real datasets. The structure of network $G(V, E)$ and objective function in (1) are explained for each dataset, followed by performance evaluation shown in plots of objective function $f(z(T))$ and disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ versus the iteration number T , where $y_i(T) = (1/T) \sum_{t=1}^T x_i(t+1)$ is the running average of $x_i(t)$ in algorithm (4) over t from 1 to T at each node i , and $z(T) = (1/m) \sum_{i=1}^m y_i(T)$ is the consensus average of $y_i(T)$ over all nodes at iteration T . For reference, we also show $f^* := f(x^*)$ in the plots of objective functions, where x^* is the optimal solution computed by MATLAB built-in linear system solvers for the synthetic decentralized least squares dataset and the real seismic datasets, and by a regular centralized gradient descent method for the synthetic robust least squares and logistic regression datasets.

4.1. Test on synthetic data. We first test on three different types of objective functions using synthetic datasets. In particular, we apply algorithm (4) to decentralized least squares, decentralized robust least squares, and decentralized logistic regression problems.

In decentralized least squares, we set the number of nodes to $m = 5$ and dimension of unknown x to $n = 5$. The radius specified for X is set to $R = 10$. For the given nodes, we generate a network by randomly turning on each of $\binom{m}{2}$ possible edges with probability 0.5 independently. For each node i , we generate a matrix A_i with $p_i = 15$ using MATLAB built-in function `randn`. We also generate a random vector $\hat{x} \in \mathbb{R}^n$ using `randn` with mean 0 and standard deviation 2. Then we simulate $b_i = A_i \hat{x} + \epsilon_i$ where ϵ_i is generated by `randn` with mean 0 and standard deviation 0.001. We set the objective function to $f_i(x) = (1/2) \|A_i x - b_i\|^2$ at node i . Therefore the Lipschitz constant of ∇f_i is $L_i = \|A_i^T A_i\|_2$, and we further set $L = \max_{1 \leq i \leq m} \{L_i\}$. The initial guess $x_i(0)$ is set to 0 for all i . For each iteration t , the delay $\tau_i(t)$ at each node i is uniformly drawn from integers 1 to B with $B = 5, 10$ and 20 . For given t , the stochastic gradient is simulated by setting $\nabla F_i(x_i(t); \xi_i(t)) = A_i^T (A_i x_i(t) - b_i) + \xi_i(t)$ where $\xi_i(t)$ is generated by `randn` with mean 0 and standard deviation σ set to 0.1 and 0.5. We run our algorithm using step size $\alpha(t) = 1/(2L + 2\eta\sqrt{t})$ with $\eta = \sqrt{[2\sigma^2 + \sqrt{\pi}G^2/\lambda^2 \log(\lambda^{-1})]/nR^2}$ which minimizes the $O(1/\sqrt{T})$ terms in the right side of (43). The objective function $f(z(T))$ and disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ versus the iteration number T are plotted in the top row of Figure 1. In the two plots, we observe that $f(z(T))$ decays to the optimal value $f^* := f(x^*)$ and disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ decays to 0 as justified by our theoretical analysis in Section 3. In general, we observe that delays with larger bound B and/or larger standard deviation σ in stochastic gradient yield slower convergence, as expected.

In the second test using synthetic dataset, we apply (4) to the decentralized robust least squares problem where the objective function is set to $f_i(x) := \sum_{j=1}^{p_i} h_i^j(x)$ with

$$(48) \quad h_i^j(x) = \begin{cases} \frac{1}{2} |(a_i^j)^T x - b_i^j|^2 & \text{if } |(a_i^j)^T x - b_i^j| \leq \delta \\ \delta \left(|(a_i^j)^T x - b_i^j| - \frac{1}{2}\delta \right) & \text{if } |(a_i^j)^T x - b_i^j| > \delta \end{cases}$$

where $(a_i^j)^T \in \mathbb{R}^n$ is the j -th row of matrix $A_i \in \mathbb{R}^{p_i \times n}$, and $b_i^j \in \mathbb{R}$ is the j -th component of $b_i \in \mathbb{R}^{p_i}$ at each node i . In this test, we simulate network $G(V, E)$ and set $A_i, b_i, m, n, R, x_i(0)$ the same way as in the decentralized least squares test above, and set $\delta = 2$ for the robust least squares. The stochastic gradient is given by $\nabla F_i(x; \xi_i(t)) = \sum_{j=1}^{p_i} \nabla h_i^j(x) + \xi_i(t)$ where $\xi_i(t)$ is generated as before with σ set to 0.1 and 0.5. Lipschitz constants L_i and L are determined as in the previous test. The settings of η and $\tau_i(t)$ remain the same as well. The objective function $f(z(T))$ and

disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ are plotted in the middle row of Figure 1. In these two plots, we observe similar convergence behavior as in the test on the decentralized least squares problem above.

The last test using a synthetic dataset is on decentralized logistic regression. In this test, we generate the network $G(V, E)$ as before but work on a slightly larger problem size where each node i possesses $A_i \in \mathbb{R}^{p_i \times n}$ with $p_i = 45$ and $n = 15$. We generate A_i the same way as before but then replace their first columns by 1. We also generate $b_i \in \{0, 1\}^{p_i}$ where each component has a random binary value. Now the objective function f_i at node i is set to

$$(49) \quad f_i(x) = \sum_{j=1}^{p_i} \left(\log[1 + \exp((a_i^j)^T x)] - b_i^j (a_i^j)^T x \right),$$

where $(a_i^j)^T \in \mathbb{R}^n$ is the j -th row of matrix $A_i \in \mathbb{R}^{p_i \times n}$, and $b_i^j \in \mathbb{R}$ is the j -th component of $b_i \in \mathbb{R}^{p_i}$. Then we perform (4) to solve this problem in the network G above. Note that $L_i \leq \|A_i^T A_i\|_2$ for all i , and we set $L = \max_{1 \leq i \leq m} \{\|A_i^T A_i\|_2\}$. The settings for the stochastic gradients, the delay $\tau_i(t)$, η , and initial value $x_i(0)$ remain the same. The objective function $f(z(T))$ and disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ are plotted in the bottom row of Figure 1, where similar convergence behavior as in the previous two tests can be observed.

4.2. Test on real data. We apply algorithm (4) to seismic tomography where the data is collected and then processed by the nodes (sensors) in a wireless sensor network. In brief, underground seismic activities (such as earthquakes) generate acoustic waves (we use P-wave here) which travel through the materials and are detected by the sensors placed on the ground. An explanatory picture of seismic tomography using a sensor network is shown in Figure 2. After data preprocessing, sensor i obtains a matrix $A_i \in \mathbb{R}^{p_i \times n}$ and a vector $b_i \in \mathbb{R}^{p_i}$, and hence an objective $f_i(x) = (1/2)\|A_i x - b_i\|^2$ for $i = 1, \dots, m$. Here $(A_i)_{kl}$, the (k, l) -th entry of matrix A_i , is the distance that the wave generated by k -th seismic activity travels through pixel l , for $k = 1, \dots, p_i$ (p_i is the total number of seismic activities) and $l = 1, \dots, n$ (n is the total number of pixels in the image), and $(b_i)_k$, the k -th component of b_i , is the total time that the wave travels from the source of k -th seismic activity to the sensor i . Then x_l , the l -th component of $x \in \mathbb{R}^n$, represents the unknown “slowness” (reciprocal of the velocity of the traveling wave) at that location (pixel) l . Once x is reconstructed from $\min_x f(x) = \sum_{i=1}^m f_i(x)$, the material (e.g., rock, sand, oil, or magma) at each pixel l can be identified by the value of x_l .

The first dataset consists of a simple and connected network G with $m = 32$ nodes where each node has 3 neighbors, and $A_i \in \mathbb{R}^{p_i \times n}$ and $b_i \in \mathbb{R}^{p_i}$ where the number of seismic events is $p_i = 512$ and the size of a 2D image x to be reconstructed is $n = 64^2 = 4096$. Since the matrix by stacking all A_i is still underdetermined, we employ an objective function with Tikhonov regularization as $f_i(x) = (1/2)(\|A_i x - b_i\|^2 + \mu\|x\|^2)$ at each node i where μ is set to 0.1. Note that more adaptive regularizers of x , such as ℓ_1 and total variation (TV) which result in a nonsmooth objective function, will be explored in future research. We apply algorithm (5) with bound B of delays set to 5, 10, and 20 and standard deviation σ of stochastic gradient to 0.5 and 0.05. The objective function $f(z(T))$ and disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ versus the iteration number T are plotted in the top row of Figure 3, where convergence of both quantities can be observed.

The second seismic dataset contains a connected network G of size $m = 50$ where each node has 3 neighbors, and matrices $A_i \in \mathbb{R}^{p_i \times n}$ and $b_i \in \mathbb{R}^{p_i}$ where $p_i = 800$ and the size of 3D image x to be reconstructed is $n = 32^3 = 32768$. We use the same objective function with Tikhonov regularization as before with $\mu = 0.01$. Other parameters are set the same as in the previous test

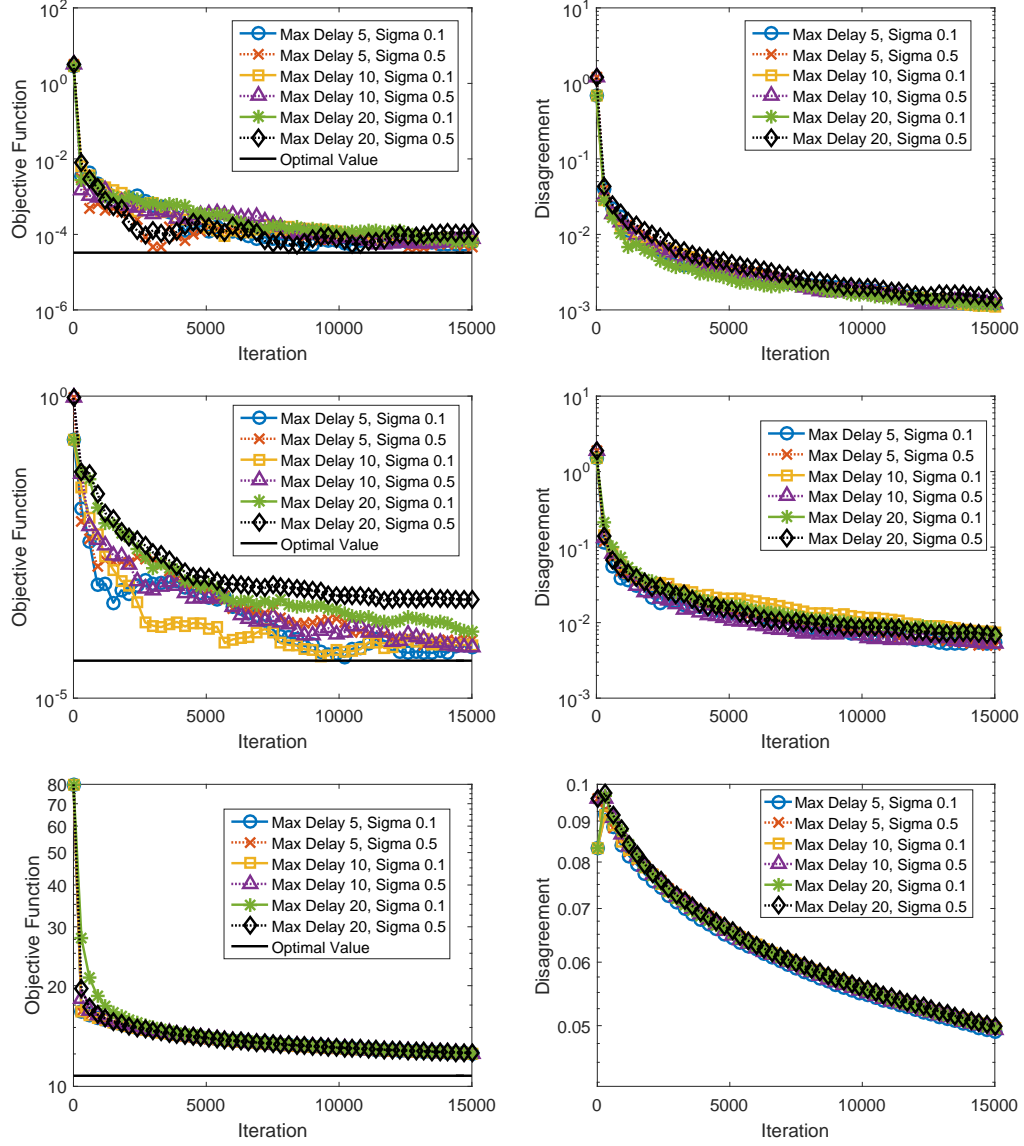
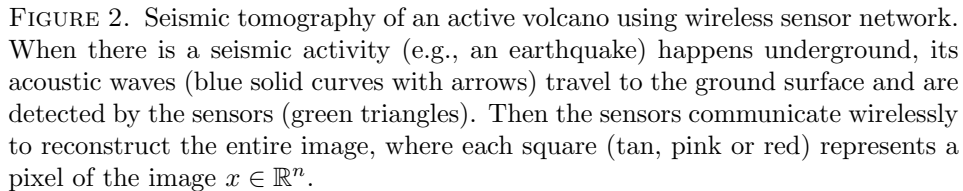


FIGURE 1. Test on synthetic decentralized least-squares (top), robust least-squares (middle), and logistic regression (bottom) for different levels of delay B and standard deviation in stochastic gradient σ . Left: objective function $f(z(T))$ versus iteration number T . Optimal value indicates $f^* := f(x^*)$. Right: disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ versus iteration number T .

on a 2D seismic image. The settings for B and σ remain the same. The objective function $f(z(T))$



The last seismic dataset consists of a connected network G of size $m = 10$ where the average node degree is 5, and matrices $A_i \in \mathbb{R}^{p_i \times n}$ and $b_i \in \mathbb{R}^{p_i}$ where $p_i = 1,816$ and the size of 3D image x to be reconstructed is $n = 160 \times 200 \times 24 = 768,000$. In this test, we employ objective $f_i(x) = (1/2)(\|A_i x - b_i\|^2 + \mu \|Dx\|^2)$ where $\mu = 0.1$ and D is the discrete gradient operator. Other parameters are set the same as in the previous two seismic datasets. The bound B of delay is set to 4, 8, and 16, and standard deviation of stochastic gradient σ is set to 1e-4 and 5e-4. The objective function $f(z(T))$ and disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ versus the iteration number T are plotted in the last row of Figure 3. The reconstructed image is displayed in the right panel of Figure 4. By comparing with the solution obtained by centralized LSQR solver (left), we can see the image is faithfully reconstructed on a decentralized network with delayed stochastic gradients.

In this paper, we analyzed the convergence of decentralized delayed stochastic gradient descent method as in (4) for solving the consensus optimization (1). The algorithm takes into consideration that the nodes in the network privately hold parts of the objective function and collaboratively solve for the consensus optimal solution of the total objective while they can only communicate with their immediate neighbors, as well as the delays of gradient information in real-world networks where the nodes cannot be fully synchronized. We show that, as long as the random delays are bounded in expectation and a proper diminishing step size policy is employed, the iterates generated by the decentralized gradient decent method converge to a consensus solution. Convergence rates of both objective and consensus were derived. Numerical results on a number of synthetic and real data were also presented for validation.

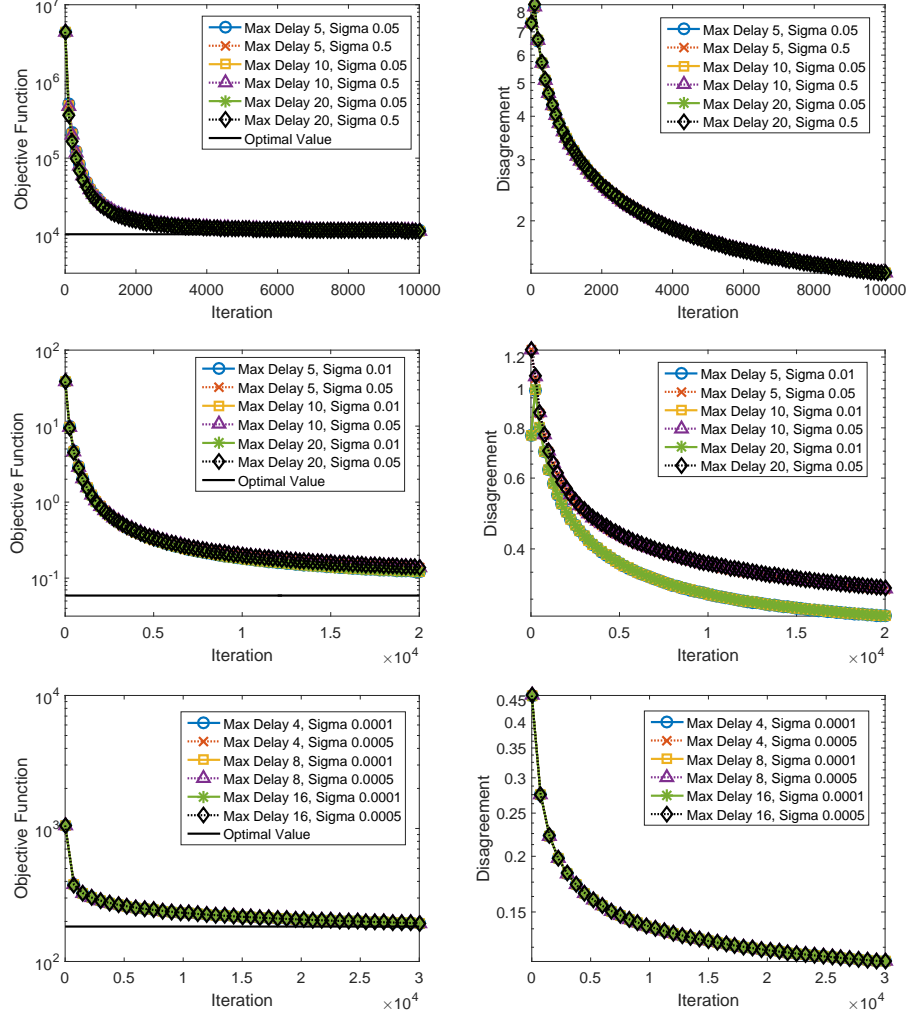


FIGURE 3. Tests on real seismic image reconstruction problems with 2D image with $n = 64^2$ (top), 3D image with $n = 32^3$ (middle), and 3D image with $n = 160 \times 200 \times 24$ (bottom) for different levels of delay B and standard deviation in stochastic gradient σ . Left: objective function $f(z(T))$ versus iteration number T . Optimal value indicates $f^* := f(x^*)$. Right: disagreement $\sum_{i=1}^m \|y_i(T) - z(T)\|^2$ versus iteration number T .

APPENDIX A. PROOF OF LEMMA 3

Proof. It suffices to show that for any fixed $R > 0$ and $X = \{x \in \mathbb{R}^m : \|x\|_\infty \leq R\}$, there is

$$(50) \quad \|(I - J)\Pi_X(x)\| \leq \|(I - J)x\|$$

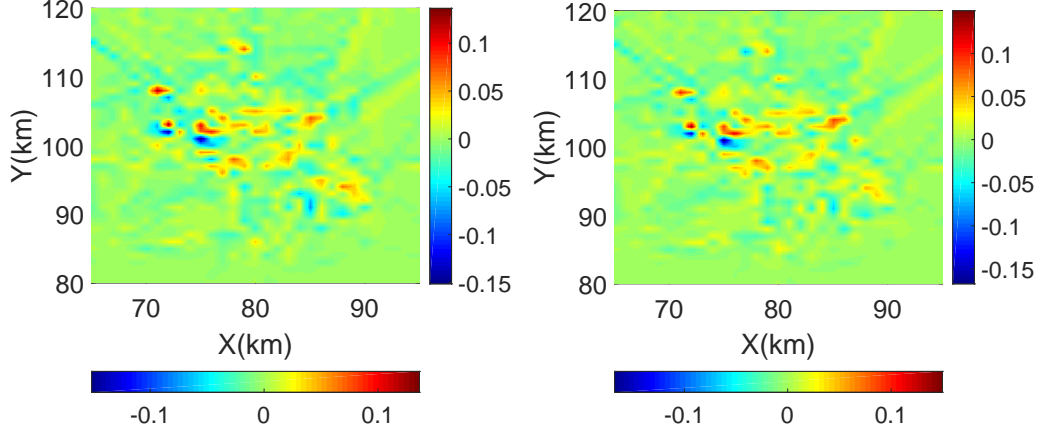


FIGURE 4. Cross section of a reconstructed 3D seismic image generated by a centralized LSQR solver (left) and decentralized algorithm with delayed stochastic gradient (4) with $B = 4$ and $\sigma = 10^{-4}$ (right).

for all $x \in \mathbb{R}^m$. Note that for $x = (x_1, x_2, \dots, x_m)^T \in \mathbb{R}^m$, there is

$$\|(I - J)x\|^2 = \sum_{i=1}^m (x_i - \bar{x})^2$$

where $\bar{x} := (1/m) \sum_{i=1}^m x_i$. We only need to show that if all $\{x_i : x_i < -R\}$ are projected to $-R$ then $\|(I - J)x\|^2$ will reduce. Without loss of generality, suppose $x_1, \dots, x_\ell < -R$ and $x_{\ell+1}, \dots, x_m \geq -R$, and let denote the means of these two groups by

$$(51) \quad \mu_1 := \frac{1}{\ell} \sum_{i=1}^{\ell} x_i < -R \quad \text{and} \quad \mu_2 := \frac{1}{m - \ell} \sum_{i=\ell+1}^m x_i \geq -R.$$

Then we have $\bar{x} = (\ell\mu_1 + (m - \ell)\mu_2)/m$, and

$$\begin{aligned}
\|(I - J)x\|^2 &= \sum_{i=1}^m (x_i - \bar{x})^2 = \sum_{i=1}^m \left(x_i - \frac{\ell\mu_1 + (m - \ell)\mu_2}{m}\right)^2 \\
&= \sum_{i=1}^{\ell} \left(x_i - \frac{\ell\mu_1 + (m - \ell)\mu_2}{m}\right)^2 + \sum_{i=\ell+1}^m \left(x_i - \frac{\ell\mu_1 + (m - \ell)\mu_2}{m}\right)^2 \\
&= \sum_{i=1}^{\ell} \left((x_i - \mu_1) + \frac{m - \ell}{m}(\mu_1 - \mu_2) \right)^2 + \sum_{i=\ell+1}^m \left((x_i - \mu_2) + \frac{\ell}{m}(\mu_2 - \mu_1) \right)^2 \\
&= \sum_{i=1}^{\ell} (x_i - \mu_1)^2 + 2 \left(\frac{m - \ell}{m} \right) (\mu_1 - \mu_2) \sum_{i=1}^{\ell} (x_i - \mu_1) + \ell \left(\frac{m - \ell}{m} \right)^2 (\mu_1 - \mu_2)^2 \\
&\quad + \sum_{i=\ell+1}^m (x_i - \mu_2)^2 + 2 \frac{\ell}{m} (\mu_2 - \mu_1) \sum_{i=\ell+1}^m (x_i - \mu_2) + (m - \ell) \left(\frac{\ell}{m} \right)^2 (\mu_2 - \mu_1)^2 \Bigg\}
\end{aligned} \tag{52}$$

After x_1, \dots, x_{ℓ} are projected to $-R$ (and $x_{\ell+1}, \dots, x_m$ remain unchanged), their mean is updated from μ_1 to $-R$ for all $i = 1, \dots, \ell$, and $\mu_2 - \mu_1 (\geq 0)$ reduces to $\mu_2 + R (\geq 0)$. Therefore, the first, third, and sixth terms in the right hand side of (52) are decreased, the second and fifth terms remain zero, and the fourth term remains unchanged. Thus $\|(I - J)x\|$ reduces after projection to $[-R, \infty)^m$. A similar argument implies that projecting $\{x_i : x_i > R\}$ to R will further reduce $\|(I - J)x\|^2$. Therefore projecting x to X , i.e., projecting to $[-R, \infty)^m$ and then $(-\infty, R]^m$, reduces $\|(I - J)x\|^2$. \square

APPENDIX B. PROOF OF LEMMA 4

Proof. First, we note that

$$\sum_{s=0}^{t-1} \alpha(s) \lambda^{t-1-s} = \alpha(0) \lambda^{t-1} + \alpha(1) \lambda^{t-2} + \sum_{s=2}^{t-1} \alpha(s) \lambda^{t-1-s} \tag{53}$$

which means that the rate is upper bounded by the last sum on the right side above since the first two tend to 0 at a linear rate $\lambda \in (0, 1)$.

Note that for all $w \in [s - 1, s]$ we have $\frac{1}{\sqrt{s}} \leq \frac{1}{\sqrt{w}}$ and $\lambda^{-s} \leq \lambda^{-(w+1)}$ since $\lambda \in (0, 1)$, and therefore

$$\alpha(s) \lambda^{t-1-s} = \frac{\lambda^{t-1-s}}{c_1 + c_2 \sqrt{s}} \leq \frac{\lambda^{t-1} \lambda^{-s}}{c_2 \sqrt{s}} \leq \frac{\lambda^{t-1} \lambda^{-(w+1)}}{c_2 \sqrt{w}} = \frac{\lambda^{t-2-w}}{c_2 \sqrt{w}}. \tag{54}$$

This inequality allows us to bound the last term on right hand side of (53) by

$$\sum_{s=2}^{t-1} \alpha(s) \lambda^{t-1-s} \leq \sum_{s=2}^{t-1} \int_{s-1}^s \frac{\lambda^{t-2-w}}{c_2 \sqrt{w}} dw = \int_1^{t-1} \frac{\lambda^{t-2-w}}{c_2 \sqrt{w}} dw = \frac{\lambda^{t-2}}{c_2} \int_1^{t-1} \frac{\lambda^{-w}}{\sqrt{w}} dw. \tag{55}$$

Now we focus on the value of integral

$$I_t := \frac{1}{2} \int_1^{t-1} \frac{\lambda^{-w}}{\sqrt{w}} dw = \int_1^{\sqrt{t-1}} \lambda^{-u^2} du \tag{56}$$

where we applied change of variables $w = u^2$. Note that we have

$$\begin{aligned}
 I_t^2 &= \int_1^{\sqrt{t-1}} \int_1^{\sqrt{t-1}} \lambda^{-(u^2+v^2)} dudv = \int_1^{\sqrt{t-1}} \int_1^{\sqrt{t-1}} e^{-(u^2+v^2) \log \lambda} dudv \\
 (57) \quad &\leq \int_0^{\sqrt{t}} \int_0^{\sqrt{t}} e^{-(u^2+v^2) \log \lambda} dudv = 2 \int_0^{\pi/4} \int_0^{\sqrt{t}/\cos \theta} e^{-\rho^2 \log \lambda} \rho d\rho d\theta \\
 &= -\frac{1}{\log \lambda} \int_0^{\pi/4} (e^{-t \log \lambda / \cos^2(\theta)} - 1) d\theta < -\frac{1}{\log \lambda} \int_0^{\pi/4} e^{-t \log \lambda / \cos^2(\theta)} d\theta
 \end{aligned}$$

where the third equality comes from changing to a polar system with the substitutions $u = \rho \cos \theta$ and $v = \rho \sin \theta$. Note that $\cos^{-2}(\theta) - (1 + 4\theta/\pi) \leq 0$ for all $\theta \in [0, \pi/4]$ since $\cos^{-2}(\theta) - 1 - 4\theta/\pi$ is convex with respect to θ and vanishes at $\theta = 0$ and $\theta = \pi/4$. Therefore

$$(58) \quad I_t^2 \leq -\frac{1}{\log \lambda} \int_0^{\pi/4} e^{-t \log \lambda (1+4\theta/\pi)} d\theta \leq \frac{\pi \lambda^{-2t}}{4t(\log \lambda)^2}.$$

Hence the sum in (55) is bounded by

$$(59) \quad \sum_{s=2}^{t-1} \alpha(s) \lambda^{t-1-s} \leq \frac{2\lambda^{t-2}}{c_2} I_t \leq \frac{2\lambda^{t-2}}{c_2} \frac{\sqrt{\pi} \lambda^{-t}}{2\sqrt{t} \log(\lambda^{-1})} = \frac{\sqrt{\pi} \lambda^{-2}}{c_2 \sqrt{t} \log(\lambda^{-1})}$$

which completes the proof. \square

REFERENCES

- [1] A. Agarwal and J. C. Duchi. Distributed delayed stochastic optimization. In *Advances in Neural Information Processing Systems*, pages 873–881, 2011.
- [2] T. C. Aysal, M. E. Yildiz, A. D. Sarwate, and A. Scaglione. Broadcast gossip algorithms for consensus. *Signal Processing, IEEE Transactions on*, 57(7):2748–2761, 2009.
- [3] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- [4] V. Cevher, S. Becker, and M. Schmidt. Convex optimization for big data: Scalable, randomized, and parallel algorithms for big data analytics. *Signal Processing Magazine, IEEE*, 31(5):32–43, 2014.
- [5] T.-H. Chang, M. Hong, W.-C. Liao, and X. Wang. Asynchronous distributed admm for large-scale optimization-part i: Algorithm and convergence analysis. *arXiv preprint arXiv:1509.02597*, 2015.
- [6] T.-H. Chang, M. Hong, and X. Wang. Multi-agent distributed optimization via inexact consensus admm. *Signal Processing, IEEE Transactions on*, 63(2):482–497, 2015.
- [7] T.-H. Chang, W.-C. Liao, M. Hong, and X. Wang. Asynchronous distributed admm for large-scale optimization-part ii: Linear convergence analysis and numerical performance. *arXiv preprint arXiv:1509.02604*, 2015.
- [8] J. C. Duchi, A. Agarwal, and M. J. Wainwright. Dual averaging for distributed optimization: convergence analysis and network scaling. *Automatic control, IEEE Transactions on*, 57(3):592–606, 2012.
- [9] H. R. Feyzmahdavian, A. Aytakin, and M. Johansson. A delayed proximal gradient method with linear convergence rate. In *Machine Learning for Signal Processing (MLSP), 2014 IEEE International Workshop on*, pages 1–6. IEEE, 2014.
- [10] P. A. Forero, A. Cano, and G. B. Giannakis. Consensus-based distributed support vector machines. *The Journal of Machine Learning Research*, 11:1663–1707, 2010.
- [11] L. Gan, U. Topcu, and S. H. Low. Optimal decentralized protocol for electric vehicle charging. *Power Systems, IEEE Transactions on*, 28(2):940–951, 2013.
- [12] F. Iutzeler, P. Bianchi, P. Ciblat, and W. Hachem. Explicit convergence rate of a distributed alternating direction method of multipliers. *arXiv preprint arXiv:1312.1085*, 2013.
- [13] F. Iutzeler, P. Ciblat, W. Hachem, and J. Jakubowicz. New broadcast based distributed averaging algorithm over wireless sensor networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*, pages 3117–3120. IEEE, 2012.

- [14] D. Jakovetic, J. M. Moura, and J. Xavier. Linear convergence rate of a class of distributed augmented lagrangian algorithms. *Automatic Control, IEEE Transactions on*, 60(4):922–936, 2015.
- [15] D. Jakovetic, J. Xavier, and J. M. Moura. Fast distributed gradient methods. *Automatic Control, IEEE Transactions on*, 59(5):1131–1146, 2014.
- [16] T. Kraska, A. Talwalkar, J. C. Duchi, R. Griffith, M. J. Franklin, and M. I. Jordan. Mlbase: A distributed machine-learning system. In *CIDR*, volume 1, pages 2–1, 2013.
- [17] J. Li, G. Chen, Z. Dong, and Z. Wu. Distributed mirror descent method for multi-agent optimization with delay. *Neurocomputing*, 2015.
- [18] M. Li, D. G. Andersen, and A. Smola. Distributed delayed proximal gradient methods. In *NIPS Workshop on Optimization for Machine Learning*, 2013.
- [19] M. Li, D. G. Andersen, A. J. Smola, and K. Yu. Communication efficient distributed machine learning with the parameter server. In *Advances in Neural Information Processing Systems*, pages 19–27, 2014.
- [20] J. Liu, S. J. Wright, C. Ré, V. Bittorf, and S. Sridhar. An asynchronous parallel stochastic coordinate descent algorithm. *The Journal of Machine Learning Research*, 16(1):285–322, 2015.
- [21] C.-H. Lo and N. Ansari. Decentralized controls and communications for autonomous distribution networks in smart grid. *Smart Grid, IEEE Transactions on*, 4(1):66–77, 2013.
- [22] A. Makhdoumi and A. Ozdaglar. Convergence rate of distributed admm over networks. *arXiv preprint arXiv:1601.00194*, 2016.
- [23] A. Mokhtari and A. Ribeiro. Decentralized double stochastic averaging gradient. *arXiv preprint arXiv:1506.04216*, 2015.
- [24] A. Nedic and A. Olshevsky. Stochastic gradient-push for strongly convex functions on time-varying directed graphs. *arXiv preprint arXiv:1406.2075*, 2014.
- [25] A. Nedic and A. Olshevsky. Distributed optimization over time-varying directed graphs. *Automatic Control, IEEE Transactions on*, 60(3):601–615, 2015.
- [26] A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *Automatic Control, IEEE Transactions on*, 54(1):48–61, 2009.
- [27] A. Nedić and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *Automatic Control, IEEE Transactions on*, 54(1):48–61, 2009.
- [28] R. Olfati-Saber and R. M. Murray. Consensus problems in networks of agents with switching topology and time-delays. *Automatic Control, IEEE Transactions on*, 49(9):1520–1533, 2004.
- [29] M. Rabbat and R. Nowak. Distributed optimization in sensor networks. In *Proceedings of the 3rd international symposium on Information processing in sensor networks*, pages 20–27. ACM, 2004.
- [30] A. Sayed. Adaptation, learning, and optimization over networks. *Foundations and Trends® in Machine Learning*, 7(4-5):311–801, 2014.
- [31] A. H. Sayed, S.-Y. Tu, and J. Chen. Online learning and adaptation over networks: More information is not necessarily better. In *Information Theory and Applications Workshop (ITA), 2013*, pages 1–8. IEEE, 2013.
- [32] O. Shamir and N. Srebro. Distributed stochastic optimization and learning. In *Communication, Control, and Computing (Allerton), 2014 52nd Annual Allerton Conference on*, pages 850–857. IEEE, 2014.
- [33] W. Shi, Q. Ling, G. Wu, and W. Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [34] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin. On the linear convergence of the admm in decentralized consensus optimization. *Signal Processing, IEEE Transactions on*, 62(7):1750–1761, 2014.
- [35] W.-Z. Song, R. Huang, M. Xu, A. Ma, B. Shirazi, and R. LaHusen. Air-dropped sensor network for real-time high-fidelity volcano monitoring. In *Proceedings of the 7th international conference on Mobile systems, applications, and services*, pages 305–318. ACM, 2009.
- [36] S. Sra, A. W. Yu, M. Li, and A. J. Smola. Adadelay: Delay adaptive distributed stochastic convex optimization. *arXiv preprint arXiv:1508.05003*, 2015.
- [37] Y.-P. Tian and C.-L. Liu. Consensus of multi-agent systems with diverse input and communication delays. *Automatic Control, IEEE Transactions on*, 53(9):2122–2128, 2008.
- [38] J. N. Tsitsiklis. Problems in decentralized decision making and computation. Technical report, DTIC Document, 1984.
- [39] H. Wang, X. Liao, T. Huang, and C. Li. Cooperative distributed optimization in multiagent networks with delays. *Systems, Man, and Cybernetics: Systems, IEEE Transactions on*, 45(2):363–369, 2015.

- [40] E. Wei and A. Ozdaglar. On the $o(1/k)$ convergence of asynchronous distributed alternating direction method of multipliers. In *Global Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE*, pages 551–554. IEEE, 2013.
- [41] L. Xiao and S. Boyd. Fast linear iterations for distributed averaging. *Systems & Control Letters*, 53(1):65–78, 2004.
- [42] D. Yuan, D. W. Ho, and S. Xu. Regularized primal-dual subgradient method for distributed constrained optimization. *IEEE Transactions on Cybernetics*, 2015.
- [43] R. Zhang and J. Kwok. Asynchronous distributed admm for consensus optimization. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1701–1709, 2014.
- [44] W. Zhang, S. Gupta, X. Lian, and J. Liu. Staleness-aware async-sgd for distributed deep learning. *arXiv preprint arXiv:1511.05950*, 2015.
- [45] L. Zhao, W.-Z. Song, L. Shi, and X. Ye. Decentralised seismic tomography computing in cyber-physical sensor systems. *Cyber-Physical Systems*, pages 1–22, 2015.
- [46] L. Zhao, W.-Z. Song, and X. Ye. Fast decentralized gradient descent method and applications to in-situ seismic tomography. In *Big Data (Big Data), 2015 IEEE International Conference on*, pages 908–917. IEEE, 2015.